



Formal Methods for Scheduling of Latency-Insensitive Designs

Julien Boucaron, Robert de Simone, Jean-Vivien Millo

► To cite this version:

Julien Boucaron, Robert de Simone, Jean-Vivien Millo. Formal Methods for Scheduling of Latency-Insensitive Designs. EURASIP Journal on Embedded Systems, 2007, 2007 (1), pp.039161. 10.1155/2007/39161 . hal-00784464

HAL Id: hal-00784464

<https://inria.hal.science/hal-00784464>

Submitted on 4 Feb 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Research Article

Formal Methods for Scheduling of Latency-Insensitive Designs

Julien Boucaron, Robert de Simone, and Jean-Vivien Millo

Aoste project-team, INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 Sophia Antipolis Cedex, France

Received 1 July 2006; Revised 23 January 2007; Accepted 11 May 2007

Recommended by Jean-Pierre Talpin

Latency-insensitive design (LID) theory was invented to deal with SoC *timing closure* issues, by allowing arbitrary fixed integer latencies on long global wires. Latencies are coped with using a *resynchronization* protocol that performs dynamic scheduling of data transportation. Functional behavior is preserved. This dynamic scheduling is implemented using specific synchronous hardware elements: *relay-stations (RS)* and *shell-wrappers (SW)*. Our first goal is to provide a formal modeling of RS and SW, that can be then formally verified. As turns out, resulting behavior is *k*-periodic, thus amenable to *static* scheduling. Our second goal is to provide formal hardware modeling here also. It initially performs *throughput equalization*, adding integer latencies wherever possible; residual cases require introduction of *fractional registers (FRs)* at specific locations. Benchmark results are presented, run on our KPASSA tool implementation.

Copyright © 2007 Julien Boucaron et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

1. INTRODUCTION

Long wire interconnect latencies induce time-closure difficulties in modern SoC designs, with propagation of signals across the die in a single clock cycle being problematic. The theory of *latency-insensitive design (LID)*, proposed originally by Carloni et al. [1, 2], offers solutions for this issue. This theory can roughly be described as such: an initial fully synchronous reference specification is first desynchronized as an asynchronous network of synchronous block components (a GALS system); it is then resynchronized, but this time with proper interconnect mechanisms allowing specified (integer-time) latencies.

Interconnects consist of fixed-sized lines of so-called *relay-stations*. These *relay-stations*, together with *shell-wrapper* around the synchronous *Pearl* IP blocks, are in charge of managing the signal value flows. With their help proper regulation of the signal traffic is performed. Computation blocks may be temporarily paused at times, either because of input signal unavailability, or because of the inability of the rest of the networks to store their outputs if they were produced. This latter issue stems from the limitation of fixed-size buffering capacity of the interconnects (*relay-station* lines).

Since their invention, *relay-stations* have been a subject of attention for a number of research groups. Extensive modeling, characterization, and analysis were provided in [3–5].

We mentioned before that the process of introducing latencies into synchronous networks introduced, at least conceptually, an intermediate asynchronous representation. This corresponds to *marked graphs* [6], a well-studied model of computation in the literature. The main property of *marked graph* is the absence of choice which matches with the absence of control in *LID*.

Marked graphs with latencies were also considered under the name of *weighted marked graphs (WMG)* [7]. We will reduce WMGs to ordinary *marked graphs* by introducing new intermediate *transportation nodes (TN)*, akin to the previous *computation nodes (CN)* but with a single input and output *link*. In fact *LID* systems can be thought of as WMGs with buffers of capacity 2 (exactly) on *link* between *computation and/or transportation nodes*. The *relay-stations* and *shell-wrappers* are an operational means to implement the corresponding flow-control and congestion avoidance mechanisms with explicit synchronous mechanisms.

The general theory of WMG provides many useful insights. In particular, it teaches us that there exists static repetitive scheduling for such computational behaviors [8, 9]. Such static *k*-periodic schedulings have been applied to software pipelining problems [10, 11], and later SoC *LID* design problems in [12]. But these solutions pay in general little attention to the form of buffering elements that are holding values in the scheduled system, and their adequacy for hardware circuit representation. We will try to provide a solution

that “perfectly” equalizes latencies over reconvergent paths, so that tokens always arrive simultaneously at the *computation node*. Sadly, this cannot always be done by inserting an integer number of latency under the form of additional transportation sections. One sometimes needs to hold back token for one step discriminatingly and sometimes does not. We provide our solution here under the form of *fractional registers (FR)*, that may hold back values according to an (input) regular pattern that fits the need for flow-control. Again we contribute explicit synchronous descriptions of such elements, with correctness properties. We also rely deeply on a syntax for schedule representation, borrowed from the theory of *N-synchronous processes* [13].

Explicit static scheduling that uses predictable synchronous elements is desirable for a number of issues. It allows a posteriori precise redimensioning of glue buffering mechanisms between local synchronous elements to allow the system to work, and this without affecting the components themselves. Finally, the extra virtual latencies introduced by equalization could be absorbed by the local computation times of *CN*, to resynthesize them under relaxed timing constraints.

We built a prototype tool for equalization of latencies and *fractional registers* insertion. It uses a number of elaborated graph-theoretical and linear-programming algorithms. We will briefly describe this implementation.

Contributions

Our first contribution is to provide a formal description of *relay-stations* and *shell-wrappers* as synchronous elements [14], something that was never done before in our knowledge (the closest effort being [15]). We introduce local correctness properties that can be easily model-checked; these generic local properties, when combined, ensure the global property of the network.

We introduce the *equalization process* to statically schedule an *LID* specification: slowing down “too fast” cycles while maintaining the original throughput of the *LID* specification. *The goal is to simplify the LID protocol.*

But rational difference of rates may still occur after *equalization process*, we solve it by adding *fractional registers (FR)*, that may hold back values according to a regular pattern that fits the need for flow-control.

We introduce a new class of *smooth* schedules that optimally minimizes the number of *FRs* used on a statically scheduled *LID* design.

Article outline

In the next section we provide some definitional and notational background on various models of computations involved in our *modeling framework*, together with an explicit representation of periodic schedules and firing instants; with this we can state historical results on *k*-periodic scheduling of *WMGs*. In Section 3, we provide the synchronous reactive representation of *relay-stations* and *shell-wrappers*, show their use in dynamic scheduling of *latency-insensitive design*, and describe several formal local correctness properties that

help with the global correctness property of the full network. Statically scheduled *LID* systems are tackled in Section 4; we describe an algorithm to build a statically scheduled *LID*, possibly adding extra virtual integer latencies and even *fractional registers*. We provide a running example to highlight potential difficulties. We also present benchmarks result of a prototype tool which implements the previous algorithms and their variations. We conclude with considerations on potential further topics.

2. MODELING FRAMEWORK

2.1. Computation nets

We start from a very general definition, describing what is common of all our models.

Definition 1 (computation network scheme). A *computation network scheme (CNS)* is a graph whose vertices are called *Computation Nodes*, and whose arcs are called *links*. We also allow arcs without a source vertex, called *input links*, or without target vertex, called *output links*.

An instance of a *CNS* is depicted on Figure 1(a).

The intention is that *computation nodes* perform computations by *consuming a data on each of its incoming links, and producing as a result a new data on each of its outgoing links*.

The occurrence of a computation thus only depends on data presence and not their actual values, so that data can be safely abstracted as *tokens*. A *CNS* is choice free.

In the sequel we will often consider the special case where the *CNS* forms a strongly connected graph, unless specified explicitly.

This simple model leaves out the most important features that are mandatory to define its operational semantics under the form of behavioral firing rules. Such features are

- (i) the *initialization* setting (where do tokens reside initially),
- (ii) the nature of *links* (combinatorial wires, simple registers, bounded or unbounded *place*, etc.),
- (iii) and the nature of *time* (synchronous, with computations firing simultaneously as soon as they can, or asynchronous, with distinct computations firing independently).

Setting up choices in these features provides distinct models of computation.

2.2. Synchronous/asynchronous versions

Definition 2. A *synchronous reactive net (S/R net)* is a *CNS* where time is synchronous: all *computation nodes* fire simultaneously. In addition *links* are either (memoryless) combinatorial wires or simple registers, and all such registers initially hold a token.

The *S/R* model conforms to synchronous digital circuits or (single-clock) synchronous reactive formalisms [16]. The network operates “at full speed”: there is always a value present in each register, so that *CN* operates at each instant.

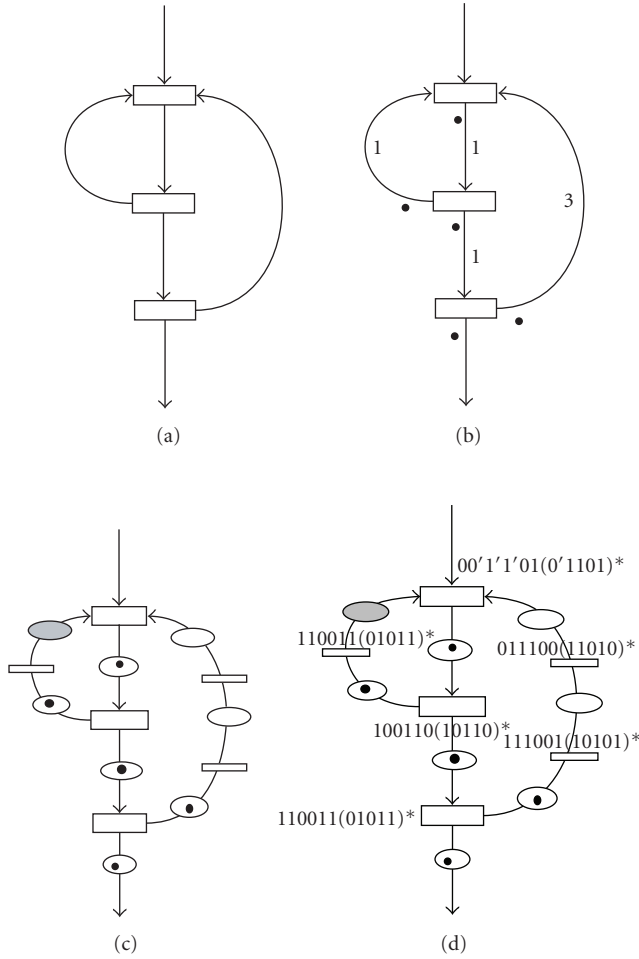


FIGURE 1: (a) An example of CNS (with rectangular *computation nodes*), (b) a corresponding WMG with latency features and token information, (c) an SMG/LID with explicit (rectangular) *transportation nodes* and (oval) *places/relay-stations*, dividing arcs according to latencies, (d) an LID with explicit schedules.

As a result, they consume all values (from registers and through wires), and replace them again with new values produced in each register. The system is *causal* if and only if there is at least one register along each cycle in the graph. Causal S/R nets are well behaved in the sense that their semantics is well founded.

Definition 3. A *marked graph* is a CNS where time is asynchronous: computations are performed independently, provided they find enough tokens in their incoming *links*; *links* have a *place* holding a number of *tokens*; in other words, *marked graphs* form a subclass of Petri Nets. The initial marking of the graph is the number of tokens held in each *place*. In addition a *marked graph* is said to be of *capacity* k if each *place* can hold no more than k tokens.

There is a simple way to encode *marked graphs* with capacity as *marked graphs* with unbounded capacity: this requires to add a reverse *link* for each existing one, which con-

tains initially a number of tokens equal to the difference between the capacity and the initial marking of the original *link*.

It was proved that a strongly connected *marked graph* is live (each computation can always be fired in the future) if and only if there is at least one token in every cycle in the graph [6]. Also, the total number of tokens in a cycle is an invariant, so strongly connected *marked graphs* are k -safe for a given capacity k .

Under proper initial conditions S/R nets and *marked graphs* behave essentially the same, with S/R systems performing all computations simultaneously “at full rate,” while similar computations are now performed independently in time in *marked graph*.

Definition 4. A *synchronous marked graph* (SMG) is a *marked graph* with an ASAP (*as soon as possible*) semantics: each *computation node* (transition) that may fire due to the availability of its input tokens immediately does so (for the current instant).

SMGs and the ASAP firing rule are underlying the works of [8, 9], even though they are not explicitly given name there.

Figure 1(c) shows a *synchronous marked graph*. Note that SMGs depart from S/R models: here all tokens are not always available.

2.3. Adding latencies and time durations

We now add latency information to indicate transportation or computation durations. These latencies will be all along constant integers (provided from “outside”).

Definition 5. A *weighted marked graph* (WMG) is a CNS with (constant integer) latency labels on *links*. This number indicates the time spent while performing the corresponding token transportation along the *link*.

We avoid computation latencies on CNs, which can be encoded as transportation latencies on *links* by splitting the actual CN into a *begin/end_CN*. Since latencies are global time durations, the relevant semantics which take them into account is necessarily ASAP. The system dynamics also imposes that one should record at any instant “how far” each token is currently in its travel. This can be modeled by an age stamp on token, or by expanding the WMG links with new *transportation nodes* (TN) to divide them into as many sections of unit latency. TNs are akin to CNs, with the particularity that they have unique source and target links. This expansion amounts to reducing WMGs to (much larger) plain SMGs. Depending on the concern, the compact or the expanded form may be preferred.

Figure 1(b) displays a *weighted marked graph* obtained by adding latencies to Figure 1(a), which can be expanded into the SMG of Figure 1(c).

For correctness matters there, still should be at least one token along each cycle in the graph, and less token on a *link* than its prescribed latency. This corresponds to the correctness required on the expanded SMG form.

Definition 6. A *latency-insensitive design (LID)* is a WMG where the expanded SMG obtained as above uses *places* of capacity 2 in between CNs and TNs.

This definition reads much differently than the original one in [2]. This comes partly from an important concern of the authors then, which is to provide a description built with basic components (named *relay-stations* and *shell-wrappers*) that can easily be implemented in hardware. Next Section 3 provides a formal representation of *relay-stations* and *shell-wrappers*, together with their properties.

Summary

CNs lead themselves quite naturally to both synchronous and asynchronous interpretations. Under some easily expected initial conditions, these variants can be shown to provide the same input/output behaviors. With explicit latencies to be considered in computation and data transportation this remains true, even if congestion mechanisms may be needed in case of bounded resources. The equivalence in the ordering of event between a synchronous circuit and an *LID* circuit is shown in [1], and equivalence between an *MG* and an *S/R* design is shown in [17].

2.4. Periodic behaviors, throughput, and explicit schedules

We now provide the definitions and classical results needed to justify the existence of static scheduling. This will be used mostly in Section 4, when we develop our formal modeling for such scheduling using again synchronous hardware elements.

Definition 7 (rate, throughput and critical cycles). Let G be a WMG graph, and C a cycle in this graph.

The *rate* R of the cycle C is equal to T/L , where T is the number of tokens in the cycle, and L is the sum of latencies of the arcs of this given cycle.

The *throughput* of the graph is defined as the minimum rate among all cycles of the graph.

A cycle is called *critical* if its rate is equal to the *throughput* of the graph.

A classical result states that, provided simple structural correctness conditions, a strongly connected WMG runs under an ultimately k -periodic schedule, with the throughput of the graph [8, 9]. We borrow notation from the theory of *N-synchronous processes* [13] to represent these notions formally, as explicit analysis and design objects.

Definition 8 (schedules, periodic words, k -periodic schedules). A *pre-schedule* for a CNS is a function $\text{Sched}: N \rightarrow w_N$ assigning an infinite binary word $w_N \in \{0, 1\}^\omega$ to every *computation node* and *transportation node* N of the graph. Node N is *activated* (or triggered, or fired, or run) at global instant i if and only if $w_N(i) = 1$, where $w(i)$ is the i th letter of word w .

A *preschedule* is a *schedule* if the allocated activity instants are in accordance with the token distribution (the

lengthy but straightforward definition is left to the reader). Furthermore, the schedule is called *ASAP* if it activates a node N whenever all its input tokens have arrived (according to the global timing).

An infinite binary word $w \in \{0, 1\}^\omega$ is called *ultimately periodic*: if it is of the form $u \cdot (v)^\omega$ where u and $v \in \{0, 1\}^*$, u represents the initialization phase, and v the periodic one.

The *length* of v is noted $|v|$ and called its *period*. The number of occurrences of 1s in v is denoted by $|v|_1$ and called its *periodicity*. The *rate* R of an ultimately periodic word w is defined as $|v|_1/|v|$.

A schedule is called k -periodic whenever for all N , w_N is a periodic word.

Thus a schedule is constructed by simulating the CNS according to its (deterministic) *ASAP* firing rule.

Furthermore, it has been shown in [9] that the length of the stationary periodic phase (called *period*) can be computed based on the structure of the graph and the (static) latencies of cycles: for a critical strongly connected component (CSCC) the length of the stationary periodic phase is the greatest common divisor (GCD) over latencies of its critical cycles. For instance assume a CSCC with 3 critical cycles having the following rates: 2/4, 4/8, 6/12, the GCD of latencies over its critical cycles is 4. For the graph, the length of its stationary periodic phase is the least common multiple (LCM) over the ones computed for each CSCC. For instance assume the previous CSCC and another one having only one critical cycle of *rate* 1/2, then the length of the stationary periodic phase of the whole graph is 2.

Figure 1(d) shows the schedules obtained on our example. If latencies were “*well balanced*” in the graph, tokens would arrive simultaneously at their consuming node; then, the schedule of any *node* should exactly be the one of its predecessor(s) shifted right by one position. However, it is not the case in general when some input tokens have to stall awaiting others. The “*difference*” (target schedule minus 1-shifted source schedule) has to be coped with by introducing specific buffering elements. This should be limited to the locations where it is truly needed. Computing the static scheduling allows to avoid adding the second register that was formerly needed everywhere in *RSs*, together with some of the backpressure scheme.

The issue arises in our running example only at the top-most *computation node*. We indicate it by prefixing some of the inactive steps (0) in its schedule by symbols: lack of input from the right input *link* (’), or from the left one (’).

3. SYNCHRONOUS TO LID: DYNAMIC SCHEDULE

In this section, we will briefly recall the theory of *latency-insensitive design*, and then focus on formal modeling with synchronous components of its main features [14].

LID theory was introduced in [1]. It relies on the fact that *links* with latency, seen as physical long wires in synchronous circuits, can be segmented into sections. Specific elements are then introduced in between sections. Such elements are called *relay-stations* (RS). They are instantiated at the oval places in Figure 1(c). Instantaneous communication

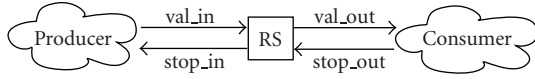


FIGURE 2: Relay-station—block diagram.

is possible inside a given section, but the values have to be buffered inside the RS before it can be propagated to the next section. The problem of computing realistic latencies from physical wire lengths was tackled in [18], where a physical synthesis floor-planner provides these figures.

Relay-stations are complemented with so-called *shell-wrappers* (SW), which compute the firing condition for their local synchronous component (called *Pearl* in LID theory). They do so from the knowledge of availability of input token and output storage slots.

3.1. Relay-stations

The signaling interface of a *relay-station* is depicted in Figure 2. The *val* signals are used to propagate tokens, the *stop* signals are used for congestion control. For symmetry here *stop_out* is an input and *stop_in* an output.

Intuitively the *relay-station* behaves as follows: when traffic is clear (no stop), each token is propagated down at the next instant from the one it was received. When a *stop_out* signal is received because of downward congestion, the RS keeps its token. But then, the previous section and the previous RS cannot be warned instantly of this congestion, and so the current RS can perfectly well receive another token at the same time it has to keep the former one. So there is a need for the RS to provide a second auxiliary register slot to store this second token. Fortunately there is no need for a third one: in the next instant the RS can propagate back a *stop_in* control information to preserve itself from receiving yet another value. Meanwhile the first token can be sent as soon as *stop_out* signals are withdrawn, and the RS remains with only one value, so that in the next step it can already allow a new one and not send its congestion control signal. Note that in this scheme there is no undue gap between the token sent.

This informal description is made formal with the description of a synchronous circuit with two registers describing the RS in Figure 3, and its corresponding syncchart [19] (in Mealy FSM style) in Figure 4. The syncchart contains the following four states.

empty when no token are currently buffered in the RS; in this state the RS simply waits for a valid input token coming, and store it in its main register that then it goes to state *half*. *stop_out* signals are ignored, and not propagated upstream, as this RS can absorb traffic.

half when it holds one token; then the RS only transmits its current, previously received token if ever does not receive an halting *stop_out* signal. If halting is requested, (*stop_out*), then it retains its token, but must also accept a potential new one coming from upstream (as it has not sent any back-pressure holding signal yet). In the second case, it becomes *full*, with the second value

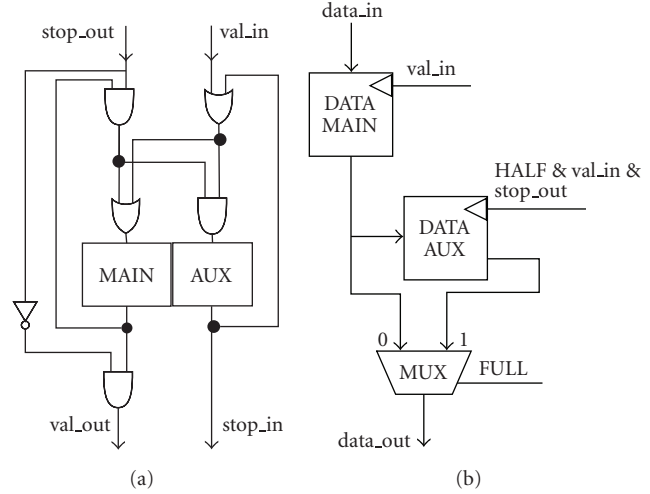


FIGURE 3: Relay-station: (a) control logic, (b) data path.

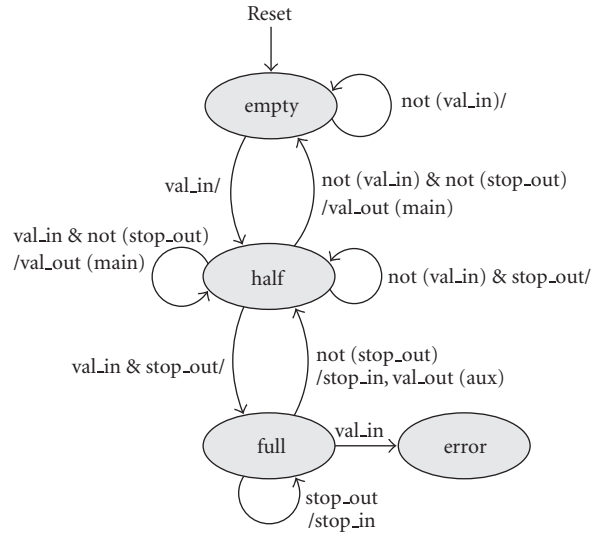


FIGURE 4: Relay-station syncchart.

occupying its “emergency” auxiliary register. If the RS can transmit (*stop_out* = *false*), it either goes back to *empty* or retrieve a new valid signal (*val_in*), remaining then in the same state. On the other hand it still makes no provision to propagate back-pressure (in the next clock cycle), as it is still unnecessary due to its own buffering capacity.

full when it contains two tokens; then it raises in any case the *stop_in* signal, propagating to the upstream section the hold-out *stop_out* signal received in the previous clock cycle. If it does not itself receive a new *stop_out*, then the line downstream was cleared enough so that it can transmit its token; otherwise it keeps it and remains halted.

error is a state which should never be reached (in an *assume/guarantee* fashion). The idea is that there should

be a general precondition stating that the environment will never send the *val_in* signal whenever the RS emits the *stop_in* signal. This should be extended to *any* combination of RS, and build up a “sequential care-set” condition on system inputs. The property is preserved as a postcondition as each RS will guarantee correspondingly that *val_out* is *not* sent when *stop_out* arrives.

NB: the notation *val_out(main)* or *val_out(aux)* means emit the signal *val_out* taking its value in the buffer, respectively, *main* or *aux*.

Correctness properties

Global correctness depends upon an assumption on the environment (see description of *error* state above). We now list a number of properties that should hold for *relay-stations*, and further *links* made of a connected line $L_n(k)$ of n successive RS elements and currently containing k values (remember that a line of n RS can store $2n$ values).

On a single RS:

- (i) $\Box \neg (stop_out \wedge val_out)$ (back-pressure control takes action immediately);
- (ii) $\Box ((stop_out \wedge X(stop_out)) \Rightarrow X(stop_in))$ (a stalled RS gets filled in two steps),

where \Box , \Diamond , \mathcal{U} , and X are the traditional *Always*, *Eventually*, *Until*, and *Next* (linear) temporal logic operators. More interesting properties can be asserted on lines of RS elements (we assume that by renaming *stop_{in,out}* and *val_{in,out}* signals form the I/O interface of the global line $L_n(k)$):

- (i) $\Box (\neg stop_out \Rightarrow \neg X^n(stop_in))$ (free slots propagate backwards);
- (ii) $\Box ((stop_out \mathcal{U} X^{(2n-k)}(true)) \Rightarrow X^{(2n-k)}(stop_in))$; (overflow);
- (iii) $(\Diamond val_in \wedge \Box (\Diamond (\neg stop_out)) \Rightarrow \Diamond val_out)$ (if traffic is not completely blocked from below from a point on, then tokens get through).

The first property is true of any line of length n , the second of any line containing initially at least k tokens, the third of any line.

We have implemented RSs and lines of RSs in the Esterel synchronous language, and model-checked combinations of these properties using *EsterelStudio*.¹

3.2. Shell-wrappers

The purpose of *shell-wrappers* is to trigger the local *computation node* exactly when tokens are available from each *input link*, and there is storage available for result in *output links*. It corresponds to a notion of *clock gating* in circuits:

the SW provides the logical clock that activates the IP component represented by the *CN*. Of course this requires that the component is physically able to run on such an irregular clock (a property called *patience* in *LID* vocabulary), but this technological aspect is transparent to our abstract modeling level. Also, it should be remembered that the *CN* is supposed to produce data on all its outputs while consuming on all its inputs in each computation step. This does not imply a combinatorial behavior, since the *CN* itself can contain internal registers of course. A more fancy framework allowing *computation latencies* in addition to our communication latencies would have to be encoded in our formalism. This can be done by “splitting” the node into *begin_CN* and *end_CN* nodes, and installing internal transportation links with desired latencies between them; if the outputs are produced with different latencies one should even split further the node description. We will not go into further details here, and keep the same abstraction level as in *LID* and *WMG* theories.

The signal interface of SWs consists of *val_in* and *stop_in* signals indexed by the number of *input links* to the SW, and of *val_out* and *stop_out* signals indexed by the number of its *output links*. There is an output *clock* signal in addition, to fire the local component. Thus, this last signal will be scheduled at the rate of local firing. Note that it is here synchronous with all the *val_out* signals when values are abstracted into tokens.

The operational behavior of the SW is depicted as a synchronous circuit in Figure 5(a), where each Input i module has to be instantiated with Figure 5(b), with its signals properly renamed, finally driving the data path in Figure 5(c). The SW is combinatorial, it takes one clock cycle to pass from RSs before the SW, through the SW and its *Pearl*, and finish into RSs in outputs of the SW. The *Pearl* is *Patient*, the state of the *Pearl* is only changed when clock (periodic or sporadic) occurs.

The SW works as follows:

- (i) the internal *Pearl's clock* and all *val_out_i* valid output signals are generated once we have all *val_in* (signal *ALL_VAL_IN* in Figure 5(a)), while *stop* is false. The internal *stop* signal itself represents the disjunction of all incoming *stop_out_j* signals from outcoming channels (signal *STOP_OUT* in Figure 5(a));
- (ii) the buffering register of a given input channel is used meanwhile as long as not all other input tokens are available (Figure 5(b));
- (iii) so, internal *Pearl's clock* is set to false whenever a backward *stop_out_j* occurs as true, or a forward *val_in_i* is false. In such case the registers already busy hold their *true* value, while others may receive a valid token “just now;”
- (iv) *stop_in_i* signals are raised towards all channels whose corresponding register was already loaded (a token was received before, and still not consumed), to warn them not to propagate any value in this clock cycle. Of course such signal cannot be sent in case the token is currently received, as it would raise a causality paradox (and a combinatorial cycle);

¹ *EsterelStudio* is a trademark of *Esterel Technologies*.

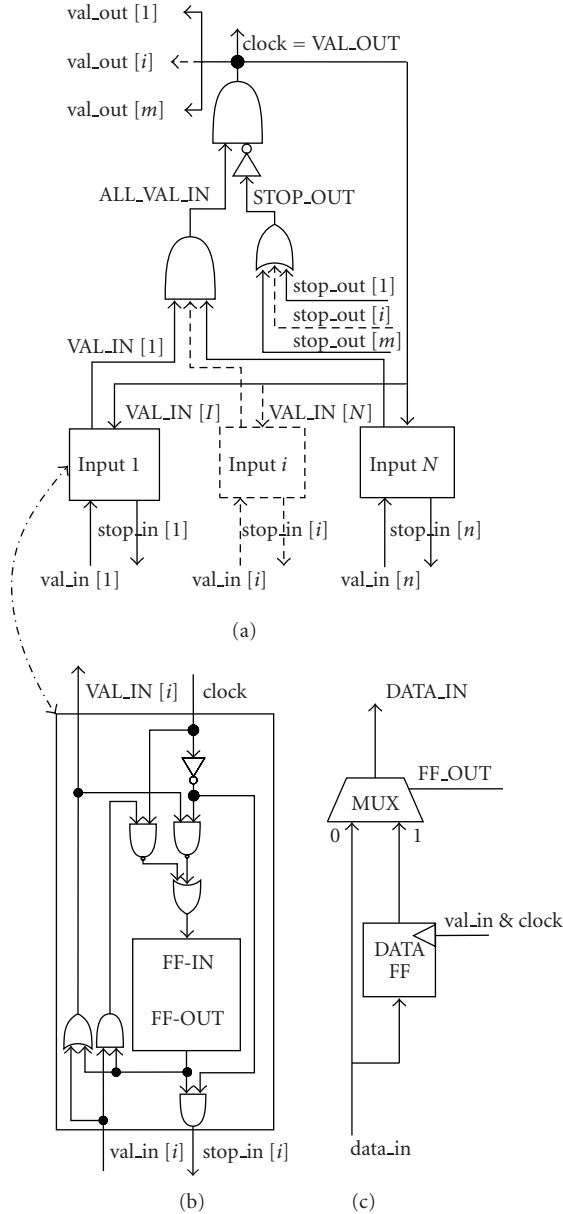


FIGURE 5: (a) Shell-wrapper circuitry, (b) input module, and (c) data path.

- (v) flip-flop registers are reset when the *Pearl's* clock is raised, as it consumes the input token. Following the previous remark, the signal $stop_in_i$ holding back the traffic in channel i is raised for these channels where the tokens have arrived before the current instant, even in this case.

Correctness properties

Again we conducted a number of model-checking experiments on SWs using ESTEREL STUDIO:

- (i) $\square ((\exists j, stop_out_j) \vee \neg clock)$ where j is an input index;

- (ii) $\square ((\exists j, stop_out_j) \Rightarrow (\forall i, \neg val_out_i))$ where j/i is an input/output index respectively;
- (iii) $\square ((\forall j, \neg stop_out_j \wedge \neg X(stop_out_j)) \Rightarrow (X(clock) \Rightarrow \exists i, X(val_in_i)))$ where j, i are input index (if the SW was not suspended at some instant by output congestion, and it triggers its pearl the next instant, then it has to be because it received a new value token on some input at this next instant).

On the other hand, most useful properties here would require syntactic sugar extensions to the logics to be easily formulated (like “a token has had to arrive on each input before or when the SW triggers its local *Pearl*,” but they can arrive in any order).

As in the case of RSs, correctness also depends on the environmental assumption that $\forall i, stop_in_i \Rightarrow \neg val_in_i$, meaning that upward components must not send a value while this part of the system is jammed.

3.3. Tool implementation

We built a prototype tool named KPASSA² to simulate and analyze an *LID* system made of a combination of previous components.

Simulation is eased by the following fact: given that the ASAP synchronous semantics of *LID* ensures determinism, for closed systems, each state has exactly one successor. So we store states that were already encountered to stop the simulation as soon as a state already visited is reached.

While we will come back to the main functions of the tool in the next section, it can be used in this context of dynamic scheduling to detect where the back-pressure control mechanisms are really been used, and which *relay-stations* actually needed their secondary register slot to preserve from traffic congestion.

4. SYNCHRONOUS TO LID: STATIC SCHEDULING

We now turn to the issue of providing static periodic schedules for *LID* systems. According to the previous philosophy governing the design of *relay-stations*, we want to provide solutions where tokens are not allowed to accumulate into *places* in large numbers. In fact we will attempt to *equalize* the flows so that tokens arrive as much as possible simultaneously at their joint *computation nodes*.

We try to achieve our goal by adding new virtual latencies on some paths that are *faster* than others. If such an ideal scheme could lead to *perfect equalization* then the second buffering slot mechanism of *relay-stations* and the back-pressure control mechanisms could be done without altogether. However, it will appear that this is not always feasible. Nevertheless, integer latency equalization provides a close approximation, and one can hope that the additional correc-

² It stands for *k-periodic ASAP Schedule Simulation and Analysis*, pronounced “*Que pasa?*”

tion can be implemented with smaller and simpler *fractional registers*.

Extra virtual latencies can often be included as computational latencies, thereby allowing the redesign of local *computation nodes* under less stringent timing budget.

As all connected graphs, general (connected) CNSs consist of directed acyclic graphs of strongly connected components. If there is at least one cycle in the net it can be shown that all cycles have to run at the rate of the slowest to avoid unbounded token accumulation. This is also true of input token consumption, and output token production rates. Before we deal with the (harder) case of strongly connected graphs that is our goal, we spend some time on the (simpler) case of acyclic graphs (with a single *input link*).

4.1. DAG case

We consider the problem of equalizing latencies in the case of directed acyclic graphs (DAGs) with a single source *computation node* (one can reduce DAGs to this sub-case if all inputs are arriving at the same instant), and no initial token is present in the DAG.

Definition 9 (DAG equalization). In this case the problem is to *equalize* the DAG such that all paths arriving to a *computation node* are having the *same latency* from inputs.

We provide a sketch of the abstract algorithm and its correction proof.

Definition 10 (critical arc). An arc is defined as *critical* if it belongs to a path of maximal latency $Max_l(N)$ from the global source *computation node* to the target *computation node* N of this arc.

Definition 11 (equalized computation node). A *computation node* N which is having only incoming *critical* arcs is defined to be an *equalized Computation Node*, that is, any path from the source to this *computation node* has the same latency $Max_l(N)$.

If a *computation node* has only one incoming arc, then this arc will be *critical* and this *computation node* will be *equalized* by definition.

The core idea of the algorithm is first to find for each *computation node* N of the graph what is its maximal latency $Max_l(N)$ and to mark incoming *critical* arcs; then the second idea is to *saturate* all *noncritical* arcs of each *computation node* of the DAG in order to obtain an *equalized* DAG.

The first part of the algorithm is done through a modified *longest-path algorithm*, marking incoming *critical* arcs for each *computation node* of the DAG and putting for each *computation node* N its maximal latency $Max_l(N)$ (as shown in Algorithm 1).

The second part of the algorithm is done as follows (see Algorithm 2). Since it may exist incoming arcs of a *computation node* N that are not *critical*, there exists an ϵ integer number that we can add such that the *noncritical* arc becomes *critical*. We can compute this integer number ϵ easily through this formula: $Max_l(N) = Max_l(N') + non_critical_arc_l + \epsilon$, where N' is the source *computation node* passing through the

```
Require: Graph is a DAG
for all ARC arc of source.getOutputArcs()
do
  NODE node  $\leftarrow$  arc.getTargetNode();
  unsigned currentLatency  $\leftarrow$ 
    arc.getLatency() + source.getLatency();
  {if the latency of this path is greater}
  if (node.getLatency()  $\leq$  currentLatency)
  then
    arc.setCritical(true);
    node.setLatency(currentLatency);
    {update arcs critical field for "node"}
    for all ARC node_arc of node.getInputArcs()
    do
      if (node_arc.getLatency() +
        node_arc.getSourceNode().getLatency() <
        currentLatency) then
        node_arc.setCritical(false);
      else
        node_arc.setCritical(true);
      end if
    end for
    {recursive call on "node" to update the whole
      sub-graph}
    recursive_longest_path(node);
  end if
end for
```

ALGORITHM 1: Procedure recursive_longest_path (NODE source).

```
Require: Graph is a DAG
for all NODE node of graph.getNodes() do
  for all ARC arc of node.getInputArcs()
  do
    if (arc.isCritical() == false) then
      unsigned maxL  $\leftarrow$  node.getLatency();
      unsigned  $\epsilon \leftarrow$  maxL
        - (arc.getLatency() +
          arc.getSourceNode().getLatency());
      arc.setLatency(arc.getLatency() +  $\epsilon$ );
      arc.setCritical(true);
    end if
  end for
end for
```

ALGORITHM 2: Procedure final_equalization (GRAPH graph).

noncritical arc and reaching the *computation node* N . Now, the *noncritical* arc through the add of ϵ is *critical*.

We apply this for all *noncritical* arcs of the *computation node* N , then the *computation node* is *equalized*.

Finally, we apply this for all *computation nodes* of the DAG, then the DAG is *equalized*.

An instance of the *unequalized*, *critical* arcs annotated and *equalized* DAG is shown in Figure 6.

Starting from the *unequalized* graph in Figure 6(a) the following holds.

The first pass of the algorithm is determining for each *computation node* its maximal latency Max_l (in circles)

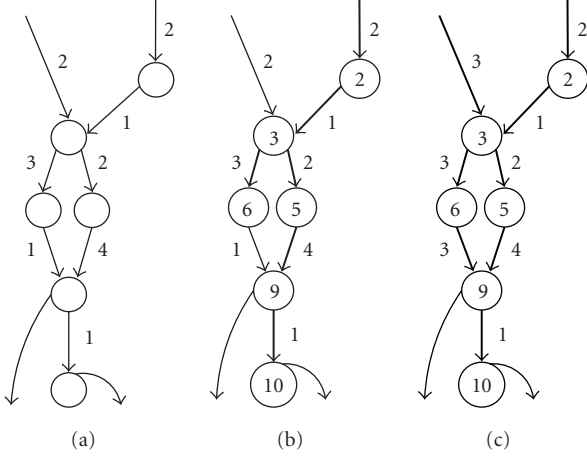


FIGURE 6: (a) *Unequalized*, (b) *critical paths annotated (large links)* and (c) *equalized DAG*.

and incoming *critical* arcs denoted using *large links* as in Figure 6(b).

The second part of the algorithm is adding “virtual” latencies (the ϵ) on *noncritical* incoming arcs, since we know the *critical* arcs coming through each *computation node* (*large links*), then we just have to add the needed amount (ϵ) in order that the *noncritical* arc is now *critical*: the sub between the value of the target *computation node*, minus the sum between the arriving *critical* arc and its source *computation node* maximal latency. For instance, consider the *computation node* holding a 9, the left branch is not *critical*, hence we are just solving $9 = 6 + 1 + \epsilon$ and $\epsilon = 2$, thus the arc will now have a latency of $3 = 1 + \epsilon$ and is so *critical* by definition. Finally, the whole graph will be fully-*critical* and thus *equalized* by definition as in Figure 6(c).

Definition 12. A *critical path* is composed only of *critical* arcs.

Theorem 1. DAG equalization algorithm is correct.

Proof. For all *computation nodes*, there is at least one *critical* arc incoming by definition; then if there is more than one incoming arc, we add the result of the sub between the maximum latency of the path passing through the so-called *critical* arc and the add between the *noncritical* arc latency and the maximum latency of the path arriving to the *computation node* where the *noncritical* arc starts. Now any arc on this given *computation node* are all *critical* and thus this *computation node* is *equalized* by definition. And this is done for any *computation node*, thus the graph is *equalized*. Since in any case we do not modify any *critical* arc, we still have the same maximum latency on *critical* paths. \square

4.2. Strongly connected case

In this case, the successive algorithmic steps involved in the process of *equalization* consist in the following:

- (1) evaluate the graph *throughput*;
- (2) insert as many additional integer latencies as possible (without changing the global *throughput*);

- (3) compute the static schedule and its initial and periodic phases;
- (4) place *fractional registers* where needed;
- (5) optimize the initialization phase (optional).

These steps can be illustrated on our example in Figure 1 as follows:

- (1) the left cycle in Figure 1(b) has *rate* $2/2 = 1$, while the (slowest) rightmost one has *rate* $3/5$. *Throughput* is thus $3/5$;
- (2) a single extra integer latency can be added to the *link* going upward in the left cycle, bringing this cycle’s *rate* to $2/3$. Adding a second one would bring the *rate* to $2/4 = 1/2$, slower than the global *throughput*. This leads to the expanded form in Figure 1(c);
- (3) the WMG is still not equalized. The actual schedules of all CN can be computed (using KPASSA, as displayed in Figure 1(d)). Inspecting closely those schedules one can notice that in all cases the schedule of a CN is the one of its predecessors shifted right by one position, *except* for the schedule of the topmost *computation node*. One can deduce from the differences in scheduling exactly when the additional buffering capacity was required, and insert dedicated *fractional registers* which delay selectively some tokens accordingly. This only happens for the initial phase for tokens arriving from the right, and periodically also for tokens arriving from the left;
- (4) it could be noticed that, by advancing only the single token at the bottom of the up going rightmost *link* for one step, one reaches immediately the periodic phase, thus saving the need for an FR element on the right cycle used only in the initial phase. Then only one FR has to be added past the regular latch register colored in grey.

We describe now the *equalization* algorithm steps in more detail.

Graph throughput evaluation

For this we enumerate all elementary cycles and compute their *rates*. While this is worst-case exponential, it is often not the case in the kind of applications encountered. An alternative would be to use well-known “minimum mean cycle problem” algorithms (see [20] for a practical evaluation of those algorithms). But the point here is that we need all those elementary cycles for setting up linear programming (LP) constraints that will allow to use efficient LP solving techniques in the next step. We are currently investigating alternative implementations in KPASSA.

Integer latency insertion

This is solved by LP techniques. Linear equation systems are built to express that all elementary cycles, with possible extra variable latencies on arcs, should now be of *rate* R , the previously computed global *throughput*. The equations are also formed while enumerating the cycles in the previous phase. An additional requirement entered to the solver can be that

the sum of added latencies be minimal (so they are inserted in a best factored fashion).

Rather than computing a rational solution and then extracting an integer approximate value for latencies, the particular shape of the equation system lends itself well to a direct *greedy* algorithm, stuffing incremental additional integer latencies into the existing systems until completion. This was confirmed by our prototype implementations.

The following example of Figure 7 shows that our integer completion does not guarantee that all elementary cycles achieve a rate very close to the extremal. But this is here because a cycle “touches” the slowest one in several distinct locations. While the global throughput is of $3/16$, given by the inner cycle, no integer latency can be added to the outside cycle to bring its rate to $1/5$ from $1/4$. Instead four fractional latencies should be added (in each arc of weight 1).

Initial- and periodic-phase schedule computations

In order to compute the explicit schedules of the initial and stationary phases we currently need to *simulate* the system's behavior. We also need to store visited state, as a termination criterion for the simulation whenever an already visited state is reached. The purpose is to build (simultaneously or in a second phase) the schedule patterns of *computation nodes*, including the quote marks (') and ('), so as to determine where residual fractional latency elements have to be inserted.

In a synchronous run each state will have only one successor, and this process stops as soon as a state already encountered is reached back. The main issue here consists in the state space representation (and its complexity). Further simplification of the state space in symbolic BDD model-checking fashion is also possible but it is out of the scope of this paper.

We are currently investigating (as “future work”) analytic techniques so as to estimate these phases without relying on this state space construction.

Fractional register insertion

In an ideally equalized system, the schedules of distinct *computation/transportation nodes* should be precisely related: the schedule of the “next” CN should be that of the “previous” CN shifted one slot right. If not, then extra *fractional registers* need to be inserted just after the regular register already set between “previous” and “next” nodes. This *FR* should delay discriminatingly some tokens (but not all).

We will introduce a formal model of our *FR* in the next subsection. The block diagram of its interfaces are displayed in Figure 8.

We conjecture that, after integer latency equalization, such elements are only required just before *computation nodes* to where cycles with different original rates reconverge. We prove in Section 4.4 that this is true under general hypothesis on smooth distribution of tokens along critical cycles. In our prototypal approach we have decided to allow them wherever the previous step indicated their need. The intention is that the combination of a regular register

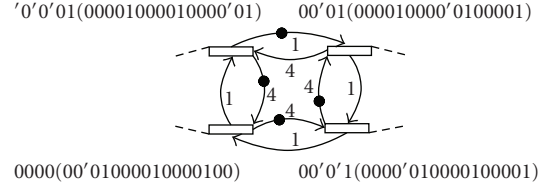


FIGURE 7: An example of WMG where no integer latency insertion can bring all the cycle rates the closest to the global throughput.

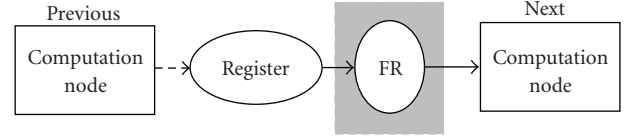


FIGURE 8: Fractional register insertion in the network.

with an additional *FR* register should roughly amount behaviorally to an *RS*, with the only difference that the backpressure control $stop_{\{in/out\}}$ signal mechanisms could be simplified due to static scheduling information computed previously.

Optimized initialization

So far we have only considered the case where all components did fire as soon as they could. Sometimes delaying some computations or transportations in the initial phase could lead faster to the stationary phase, or even to a distinct stationary phase that may behave more smoothly as to its scheduling. Consider in the example of Figure 1(c) the possibility of firing the lower-right *transportation node* alone (the one on the backward up arc) in a first step. This modification allows the graph to reach immediately the stationary phase (in its last stage of iteration).

Initialization phases may require a lot of buffering resources temporarily that will not be used anymore in the stationary phase. Providing short and buffer-efficient initialization sequences becomes a challenge. One needs to solve two questions: first, how to generate efficiently states reachable in an *asynchronous* fashion (instead of the deterministic *ASAP* single successor state); second, how to discover very early that a state may be part of a periodic regime. These issues are still open. We are currently experimenting with *KPASSA* on efficient representation of *asynchronous* firings and resulting state spaces.

Remark 1. When applying these successive transformation and analysis steps, which may look quite complex, it is predictable that simple subcases often arise, due to the well-chosen numbers provided by the designer. Exact integer equalization is such a case. The case when fractional adjustments only occur at reconvergence to critical paths are also noticeable. We built a prototype implementation of the approach, which indicates that these specific cases are indeed often met in practice.

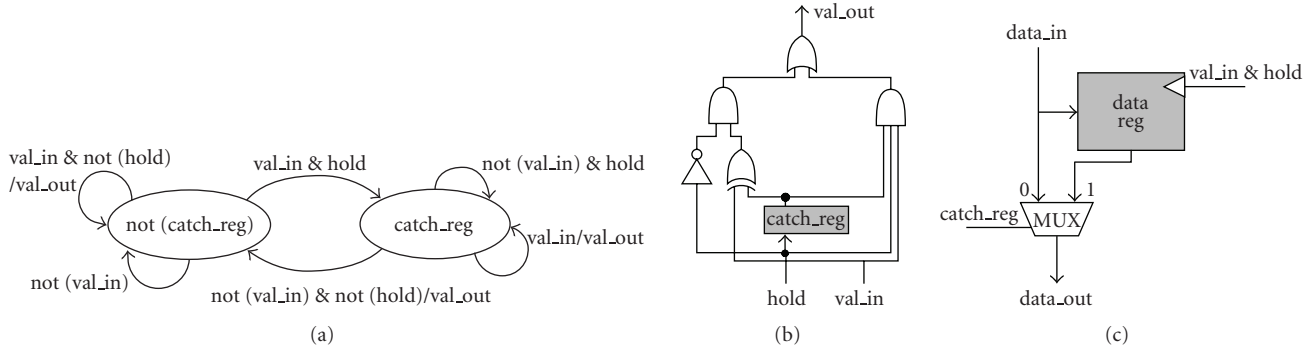


FIGURE 9: (a) The syncchart, (b) the interface block-diagram of the FR, and (c) the datapath.

4.3. Fractional register element (FR)

We now formally describe the specific *FR*, both as a synchronous circuit in Figure 9(b) and as a corresponding syncchart (in Mealy FSM style) in Figure 9(a).

The *FR* interface consists of two input wires *val.in* and *hold*, and one output wire *val.out*. Its internal state consists of a register *catch_reg*. The register will be used to “kidnap” the valid data (and its value in a real setting) for one clock cycle whenever *hold* holds. We note *pre(catch_reg)* the (boolean) value of the register computed at the previous clock cycle. It indicates whether the slot is currently occupied or free.

It is possible that the same data is held several instants in a row. But meanwhile there should be no new data arriving, as the *FR* can store only one value; otherwise this would cause a conflict.

It is also possible that a full sequence of consecutive data are held back one instant each in a burst fashion. But then each data/value should leave the element in the very next instant to be consumed by the subsequent *computation node*; otherwise this would also cause a conflict.

Stated formally, when $hold \wedge pre(catch_reg)$ holds then either *val.in* holds, in which case the new data enters and the current one leaves (by scheduling consistency the *computation node* that consumes it should then be active), or *val.in* does not hold, in which case the current data remains (and, again by scheduling consistency, then the *computation node* should be inactive). Furthermore the two extra conditions are requested:

$[hold \Rightarrow (val_in \vee pre(catch_reg)):]$ if nothing can be held, the scheduling does not attempt to;

$[(val_in \wedge pre(catch_reg)) \Rightarrow hold:]$ otherwise the two pieces of data could cross the element and be output simultaneously.

The *FR* behavior amounts to the two equations:

$[catch_reg = hold:]$ the register slot is used only when the scheduling demands;

$[val_out = val_out_1 \vee val_out_2:]$

- (i) $val_out_1 = val_in \oplus pre(catch_reg) \wedge \neg hold;$
- (ii) $val_out_2 = val_in \wedge pre(catch_reg) \wedge hold.$

either a new value directly falls across, or an old one is chased by a new one being held in its *place*.

Our main design problem is now to generate *hold* signals exactly when needed. Its schedule should be the difference between the schedule of its source (*computation or transportation node*) shifted by one instant, and the schedule of its target node; indeed, a token must be held when the target node does not fire while the source *CN* did fire to produce a token last instant, or if the token was already held at last instant.

Consider again Figure 8, we will name *w* the schedule of the *previous* source *CN*, and *w'* the schedule of the *next* target *CN*. After the regular register delay the data are produced to the *FR* entry on schedule $0.w$ (shifted one slot/instant right). The *fractional register* should hold the data exactly when the *k*th active step at this entry is not the *k*th activity step at its target *CN* that must consume it. In other words, the *FR* resynchronize its input and output, which cannot be away more than one activity step. This last property is true as the schedules were computed using the *LID* approach with *relay-stations*, which do not allow more than one extra token in addition to the regular one on each arc between *computation or transportation nodes*.

Stated formally, this property becomes: $hold(n) = 1$ if and only if $|0 \cdot w_n|_1 \neq (|w'_n|_1 - |w'_0|_1)$. It says that at a given instant *n* we should kidnap a value if the number of occurrences of 1 up to instant *n* on the previous *CN* is different than the number of occurrences of 1 on the next *computation node*. More precisely, the $-|w'_0|_1$ term takes care of a possible initial activity at the target *CN*, not caused by the propagation of tokens from the source *CN*, that would have to be removed.

Figure 10 shows a possible implementation computing *hold* from signals that would explicitly provide the target and source schedules as inputs.

Correctness properties

It can be formally proved that, under proper assumptions, a full RS is sequentially equivalent to a system made of a regular register followed by a fractional one, with the respective *stop_out* and *hold* signals equated (as in Figure 11). The exact assumption is that a *stop_out/hold* signal is never received when the systems considered are already full (both

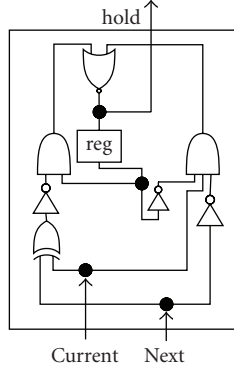


FIGURE 10: Hold implementation.

registers occupied in each case). Providing this assumption to a model-checker is cumbersome, as it deals with internal states. It can thus be replaced by the fact that never in history there are more than one *val_in* signal received in excess of the *val_out* signals sent. This can easily be encoded by a synchronous observer.

In essence the previous property states that the two systems are equivalent *safe* for the emission of *stop_in* on a full RS. This emission can also be shown to be simulated by inserting the previous *HOLD* component with proper inputs. Of course, this *does not* mean that the implementation will use such a dynamic *HOLD* pattern, but that simulating its effect (because the static scheduling instructs us of when to generate the signal) would make things equal to the former RS case.

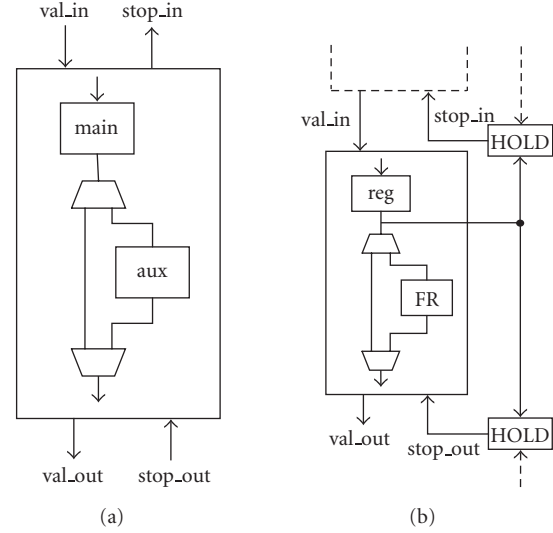
4.4. Issues of optimal FR allocation

As already mentioned in the case of an SCC we still do not have a proof that in the stationary phase it is enough to include such elements at the entry points of *computation nodes* only, so they can be installed in place of more *relay-stations* also. Furthermore, it is easy to find initialization phases where tokens in excess will accumulate at any locations, before the rate of (the) slowest cycle(s) distributes them in a smoother, evenly distributed pattern. Still we have several hints that partially deal with the issue. It should be remembered here that, even without the result, we can equalize latencies (it just needs adding more FRs).

Definition 13 (smoothness). A schedule is called *smooth* if the sequences of successive 0 (inactive) instants the difference in length between sequences of consecutive 0s cannot differ by more than 1. The schedule (1001)* is *not* smooth since they are two consecutive 0 between the first and second occurrences of 1, while there is none between the second and the third.

Conjecture 1. If all computation node schedules are smooth, rates can be equalized using FR only at computation node entry points.

Counter example 1. We originally thought that Conjecture 1 should be sufficient, but the counter example of Figure 12



SHIFT(): the data in register “reg” goes in the “FR.”
It is an internal function.

FIGURE 11: Equivalence of RS and FR roles.

was found. Assume a simple graph formed with two cycles sharing one CN. The first critical cycle has 7 tokens and 11 latencies, the second one has 5 tokens and 7 latencies. There exists a stationary phase where the schedule of all CNs is smooth (it is [101010111] or any rotation of this word) but we need two successive FRs on the noncritical cycle because only one FR should overflow.

The reason of this failure is that the definition of smoothness is not restrictive enough. In the schedule of the counterexample Figure 12, the pattern 10 is repeated 3 times at the beginning and we have 3 occurrences of 1 (which are not followed by any 0) at the end. 0 and 1 are not spread regularly

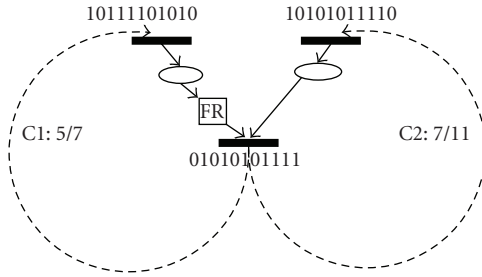


FIGURE 12: Counter example of Conjecture 1. The FR overflow at instant 7.

enough in the schedule. However, if the schedule of the CN become (01011011011), we now need only one *FR* on the noncritical cycle.

We propose a *new* definition.

Definition 14 (extended smoothness). A schedule w is said to be *extended smooth* if any subword, with a length l , contains either n bits at 1 or $n + 1$ bits at 1, where n is equal to $\lfloor l * |w|_1 / |w| \rfloor$, $|w|_1$ is the number of occurrences of 1 in w and $|w|$ is the length of w .

4.5. Tool implementation

Our KPASSA tool implements the various algorithmic stages described above. Given that we could not yet prove that *FR*s were only required at specific locations, the tool is ready to insert some anywhere. KPASSA computes and displays the system throughput, showing critical cycles and the locations of choice for extra integer latency insertions in noncritical cycles. It then computes an explicit schedule for each *computation and transportation node* (in the future it could be helpful to display only the important ones), and provides locations for *fractional registers* insertion. It also provides log information on the numbers of elements added, and whether perfect integer equalization was achieved in the early steps.

In the future, we plan to experiment with algorithms for finding efficient asynchronous transitory initial phases that may reach the stationary periodic regime faster than with the current *ASAP* synchronous firing rule.

Figure 13 displays a screen copy of KPASSA on a case study drawn from [3]. Using the original latency specifications our tool found a static schedule using less resources than the former implementation based on *relay-stations* and dynamic back-pressure mechanisms. And now the activation periods of components are fully predictable.

5. EXPERIMENTS ON CASE STUDIES

Tables 1 and 2 display benchmark results obtained with KPASSA on a number of case studies. The first examples were built from [3] for MPEG2 video encoder and from existing and publicly available models of structural IP block diagrams (IP MegaStore of Altera). But the latency figures were suggested by our industrial partners of PACA CIM initiative. In [18] the authors use a public-domain floorplanner to syn-

thesize approximate latency figures, based on wire lengths induced by the placement of IPs. The last two examples are based on graph shapes and latency distribution that are a priori adverse to the approach (without being formerly worst-cases).

Table 1 provides features of size that are relevant to the algorithmic complexity. Table 2 reports the results obtained, about whether perfect equalization holds, the number of *fractional registers* required in the initial and periodic phases (note that some *FR* elements may still be needed for the initial part even in perfectly equalized cases), the number of integer latencies added, and time and space performances.

The current implementation of the tool is not yet optimized for complexity in time and space, until now this is not yet important. The graph state encoding is naive, and algorithms are not optimal.

KPASSA is a formal tool that is able to compute effectively the length of initialization and periodic patterns to compute an upper-bound of the number of resources used for the implementation. The tool provides huge preliminary implementations for the static-scheduled LID, but it let us experiment new ideas to optimize those implementations.

In addition to the results shown in Tables 1 and 2, KPASSA also provides synthetic information on the criticality of nodes: cycles can be ordered by their rates, and then nodes by the slowest rate of a cycle it belongs to. Then the nodes are painted from red “(Hotspot)” to blue “(Coldspot)” accordingly. This visual information is particularly useful before *equalization*.

6. FURTHER TOPICS

Concerning the static scheduling, a number of important topics are left open for further theoretical developments as follows.

- (i) Relaxing the firing rule: so far the theory developed here only considers the case where local synchronous components all consume and produce token on all input and output channels in each computation step, and where they all run on the same clock. In this favorable case functional determinacy and confluence are guaranteed, with latencies only impacting the relative ordering of behaviors. So it can be proved that the relaxed-synchronous version produces the same output streams from the same input streams as the fully synchronous specification (indeed the rank of a token in a stream corresponds to its time in the synchronous model, thereby reconstructing the structure of successive instants). Several papers considered extensions in the context of GALS systems, but then ignored the issue of functional correspondence with an initial well-clocked specification, which is our important correctness criterion. This relaxation may help minimize some metrics:

- (a) we certainly would like to establish that *FR* are needed only at *computation nodes*, minimizing their number rather intuitively;

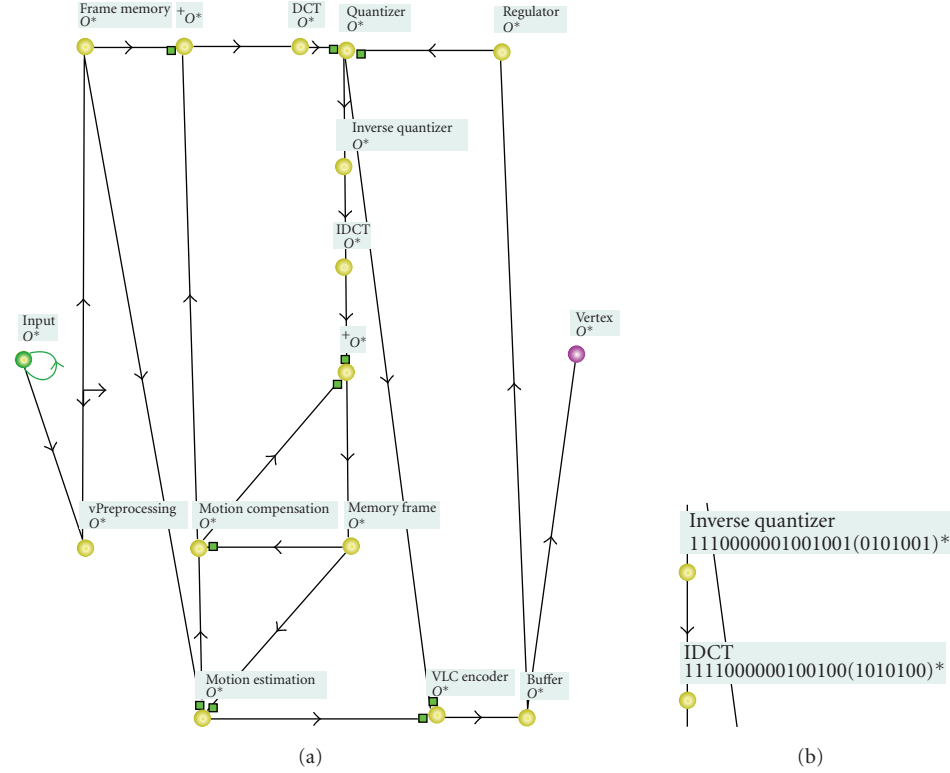


FIGURE 13: An example simulation result (MPEG2 Encoder) with KPASSA. In (a), the graph; in (b), the displayed schedules for two vertices.

TABLE 1: Example sizes before equalization.

	No. of nodes	No. of cycles	No. of critical cycles	Max cycle latency	Throughput
MPEG2 video encoder	16	7	3	21	3/7
Encoder multistandard ADPCM	12	23	23	14	1/2
H264/AVC encoder	20	12	3	27	4/9
29116a 16 bits CAST MicroCPU	11	7	3	35	3/35
Abstract stress cycles	40	2295	1	1054	4/29
Abstract stress nodes	175	3784	1	1902	4/29

(b) discovering short and efficient (minimizing number of *FR*) initial phases is also an important issue here;

(c) the distribution of integer latencies over the arcs could attempt to minimize (on average) the number of *computation nodes* that are active altogether. In other words, transportation latencies should be balanced so that computations alternate in time whenever possible. The goal is here to avoid “*hot spots*” that is to say flatten the power peaks. It could be achieved by some sort of retiming/recycling techniques and schedules exploration still using a relaxed firing rule.

(ii) *Marked graphs* do not allow control-flow (and control *modes*). The reason is, in general case such as full Petri Nets, it can no longer be asserted that tokens are consumed and produced at the same rate. But explicit “*branch schedules*” could probably help regulate

the branching control parts by a way similar to that by which they control the flow rate.

Finally, the goal would be to define a general GALS modeling framework, where GALS components could be put in GALS networks (to this day the framework is not compositional in the sense that local components need to be synchronous). A system would consist again of computation and interconnect communication blocks, this time each with appropriate triggering clocks, and of a scheduler providing the subclocks computation mechanism, based on their outer main clock and several signals carrying information on control flow.

Summary

In this article we first introduced full formal models of relay stations and Shell Wrappers, the basic components for the theory of latency-insensitive design. Altogether they allow to

TABLE 2: Equalization performances and results (run on P4 3.4 GHz, 1 GB RAM, Linux 2.6, and JDK 1.5).

	Perfect eqn.	No. of FR init./periodic	No. of added latencies	Time	Memory
MPEG2 video encoder	N	9/5	18	< 1 s	~11 MB
Encoder multistandard ADPCM	Y	24/0	91	< 1 s	~11 MB
H264/AVC encoder	N	18/11	0	~1 s	~11 MB
29116a 16 bits CAST MicroCPU	Y	0/0	0	~1 s	~11 MB
Abstract stress cycles	N	55/24	1577	~17 s	~16 MB
Abstract stress nodes	N	59/23	2688	~4 min	~43 MB

build a dynamic scheduling scheme which stalls traveling values in case of congestion ahead. We established a number of correctness properties holding between (lines of) RSs and SWs.

Then, using former results from scheduling theory we recognized the existence of static periodic schedules for networks with fixed constant latencies. We tried to use these results to compute and optimize the allocation of buffering resources to the system. By *equalization* we obtain location where a full extra latency is *always* mandatory (these virtual latencies can be later absorbed in the redesign of more relaxed IP components). Fractional latencies still need to be inserted to provide *perfect* equalization of throughputs. By simulation we compute the exact schedules of computation nodes, and deduce the locations of fractional register assignments to support that. We conjectured that under simple “smoothness” assumptions on the token values distribution along graph cycles, the FR elements could be inserted in an optimized fashion. We also proved properties on FR implementation, and its relation to RSs.

Finally, we described a prototype implementation of the techniques used to compute schedules and allocate integer and fractional latencies to a system, together with preliminary benchmarks on several case studies.

ACKNOWLEDGMENTS

This work was partially supported by ST Microelectronics and Texas Instruments grants in the context of the French regional PACA CIM initiative.

REFERENCES

- [1] L. P. Carloni, K. L. McMillan, and A. L. Sangiovanni-Vincentelli, “Theory of latency-insensitive design,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 20, no. 9, pp. 1059–1076, 2001.
- [2] L. P. Carloni, K. L. McMillan, A. Saldanha, and A. L. Sangiovanni-Vincentelli, “A methodology for correct-by-construction latency insensitive design,” in *Proceedings of the IEEE/ACM International Conference on Computer-Aided Design (ICCAD ’99)*, pp. 309–315, San Jose, Calif, USA, November 1999.
- [3] L. P. Carloni and A. L. Sangiovanni-Vincentelli, “Performance analysis and optimization of latency insensitive systems,” in *Proceedings of the 37th Conference on Design automation (DAC ’00)*, pp. 361–367, Los Angeles, Calif, USA, June 2000.
- [4] T. Chelcea and S. M. Nowick, “Robust interfaces for mixed-timing systems with application to latency-insensitive protocols,” in *Proceedings of the 38th conference on Design automation (DAC ’01)*, pp. 21–26, Las Vegas, Nev, USA, June 2001.
- [5] A. Chakraborty and M. R. Greenstreet, “A minimalist source-synchronous interface,” in *Proceedings of the 15th Annual IEEE International ASIC/SOC Conference*, pp. 443–447, Rochester, NY, USA, September 2002.
- [6] F. Commoner, A. W. Holt, S. Even, and A. Pnueli, “Marked directed graphs,” *Journal of Computer and System Sciences*, vol. 5, no. 5, pp. 511–523, 1971.
- [7] C. Ramchandani, *Analysis of asynchronous concurrent systems by timed Petri nets*, Ph.D. thesis, MIT, Cambridge, Mass, USA, September 1973.
- [8] J. Carlier and P. Chrétienne, *Problème d’ordonnancement: modélisation, complexité, algorithmes*, Masson, Paris, France, 1988.
- [9] F. Baccelli, G. Cohen, G. J. Olsder, and J.-P. Quadrat, *Synchronization and Linearity: An Algebra for Discrete Event Systems*, John Wiley & Sons, New York, NY, USA, 1992.
- [10] V. van Dongen, G. R. Gao, and Q. Ning, “A polynomial time method for optimal software pipelining,” in *Proceedings of the 2nd Joint International Conference on Vector and Parallel Processing (CONPAR ’92)*, pp. 613–624, Springer, Lyon, France, September 1992.
- [11] F.-R. Boyer, E. M. Aboulhamid, Y. Savaria, and M. Boyer, “Optimal design of synchronous circuits using software pipelining techniques,” in *Proceedings of IEEE International Conference on Computer Design (ICCD ’98)*, pp. 62–67, Austin, Tex, USA, October 1998.
- [12] M. R. Casu and L. Macchiarulo, “A new approach to latency insensitive design,” in *Proceedings of the 41st Annual Conference on Design Automation (DAC ’04)*, pp. 576–581, ACM Press, San Diego, Calif, USA, June 2004.
- [13] A. Cohen, M. Duranton, C. Eisenbeis, C. Pagetti, F. Plateau, and M. Pouzet, “N-synchronous Kahn networks: a relaxed model of synchrony for real-time systems,” in *Proceedings of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL ’06)*, pp. 180–193, ACM Press, Charleston, South Carolina, USA, January 2006.
- [14] J. Boucaron, J.-V. Millo, and R. de Simone, “Another glance at relay stations in latency-insensitive design,” *Electronic Notes in Theoretical Computer Science*, vol. 146, no. 2, pp. 41–59, 2006.
- [15] M. R. Casu and L. Macchiarulo, “A detailed implementation of latency insensitive protocols,” in *Proceedings of Formal Methods for Globally Asynchronous Locally Synchronous Architectures*, pp. 94–103, Pisa, Italy, September 2003.
- [16] A. Benveniste, P. Caspi, S. A. Edwards, N. Halbwachs, P. Le Guernic, and R. de Simone, “The synchronous languages 12 years later,” *Proceedings of the IEEE*, vol. 91, no. 1, pp. 64–83, 2003.

- [17] A. V. Yakovlev, A. M. Koelmans, and L. Lavagno, "High-level modeling and design of asynchronous interface logic," *IEEE Design and Test of Computers*, vol. 12, no. 1, pp. 32–40, 1995.
- [18] M. R. Casu and L. Macchiarulo, "Floorplanning for throughput," in *Proceedings of the International Symposium on Physical Design (ISPD '04)*, pp. 62–69, ACM Press, Phoenix, Ariz, USA, April 2004.
- [19] C. André, "Representation and analysis of reactive behaviors: a synchronous approach," in *Proceedings of the IMAC Multiconference on Computational Engineering in Systems Applications (CESA '96)*, pp. 19–29, Lille, France, July 1996.
- [20] A. Dasdan, "Experimental analysis of the fastest optimum cycle ratio and mean algorithms," *ACM Transactions on Design Automation of Electronic Systems*, vol. 9, no. 4, pp. 385–418, 2004.